

What is claim is:

1. An automatic method of classifying molecules having similar biologic function comprising the steps of: a) creating a hierarchical organization of said molecules in a database, wherein groups of clusters are identified using local consideration resulting in related clusters; b) determining the position of a selected molecule based on the hierarchical organization of step a, whereby selected molecules of similar biologic function are classified.
2. An automatic method of classifying molecules having similar biologic function comprising the steps of:
 - a) creating a hierarchical organization of the molecules in a database comprising the steps of:
 - i calculating pairwise similarities between said molecules in said database by combining at least one standard measure of similarity, resulting in a first set of expectation values of similarity;
 - ii analyzing said first set of expectation values of similarity so as to obtain a second set of expectation values of similarity, wherein said molecules of the second set have a high degree of similarity;

iii merging said resulting molecules of step ii) so as to form clusters, wherein only molecules with an expectation value below a first restricted threshold are merged;

iv identifying groups of clusters using local consideration resulting in related clusters;

v determining the relationship between related clusters in said groups;

vi analyzing said groups of the related clusters of step v), thereby creating a hierarchical organization of said molecules;

b) determining the position of a selected molecule based on the hierarchical organization of step a, comprising the steps of:

i identifying pairwise similarities between said selected molecule to the molecules in said database by combining at least one standard measure of similarity, resulting in a third set of expectation values of similarity;

ii identifying geometric averaging of said selected molecule to each of said resulting clusters in said hierarchy, resulting in a fourth set of expectation values of similarity;

iii identifying related clusters from said hierarchy of step a, having a geometric averaging with said selected molecule below a second threshold, thereby classifying molecules having similar biologic function.

3. The method of claim 2, wherein each related cluster resulting from step iv is further analyzed, wherein if said geometric averaging of said related cluster to a second related cluster is below said first threshold a connection is established between two clusters.
4. The method of claim 3, wherein said cluster is analyzed by calculating the geometric averaging between molecules in one cluster to molecules in the other cluster.
5. The method of claim 2, wherein the method of determining the relationship between said related clusters in step v further comprises the steps of applying a global test on said connected clusters for identifying nuclei of strong relationships within said groups of clusters, wherein each cluster is checked against its nearest cluster by using said geometric averaging, wherein if said geometric averaging is below said first threshold the two clusters merged.
6. The method of claim 2, further comprises repeating steps iv, v while raising the threshold in a step wise manner until a third threshold is achieved.
7. The method according to claim 2, wherein said groups containing the clusters of every two molecules from different clusters with expectation value below the threshold.
8. The method according to claim 2, wherein the method further comprises applying three standard measures of similarity.
9. The method according to claim 8, wherein said standard measures of similarity are Smith Waterman (SW), FASTA, and BLAST.

10. The method according to claim 9, further comprising the step of applying a numerical normalization to the SW, the FASTA and the BLAST.
11. The method according to claim 10, wherein the method step a) ii) is analyzed using thresholds 0.1, 0.1 and 10^{-3} for SW, FASTA and BLAST respectively wherein a value of pairwise similarities between two protein sequences is maintained if either SW or FASTA yield an expectation value ≤ 0.1 or BLAST's expectation value is $\leq 10^{-3}$.
12. The method according to claim 11, further comprising filtering the results of BLAST thereby excluding low complexity segments, using the SEG program; or setting a more stringent threshold for BLAST at 10^{-6} if filtering decreased the number of related sequences by half.
13. The method according to claim 12, further comprising maintaining a value of an expectation value of similarity between two protein when all three methods yield an expectation value ≤ 1 .
14. The method of claim 1, wherein the molecules are nucleic acids.
15. The method of claim 1, wherein the nucleic acids are DNA or RNA.
16. The method of claim 1, wherein the molecules are polynucleotides.
17. The method of claim 1, wherein the molecules are proteins.
18. The method of claim 1, wherein the molecules are peptides.
19. The method of claim 1, wherein the molecules are glycoproteins.
20. The method of claim 1, wherein the molecules are complex sugars.
21. The method of claim 1, wherein the molecules are immunoglobulins.

22. The method of claim 1, wherein the molecules are chemicals.
23. The method of claim 1, wherein said database is the Swissprot database.
24. The method of claim 2, wherein said first threshold is E^{-100} .
25. The method of claim 2, wherein said second threshold is E^{10} .
26. The method of claim 2, wherein said third threshold is E^0 .
27. A system for identifying biological families of molecules, said system comprising :
- a) a storage device containing a first database identifying molecules and their corresponding hierarchy;
 - b) a processor connected to said storage device and in communication with a second database containing molecule sequences;
- an input device connected to said processor for inputting a selected molecule to be analyzed, wherein said processor is programmed to automatically classify molecules having similar biologic function.

09601278-022201